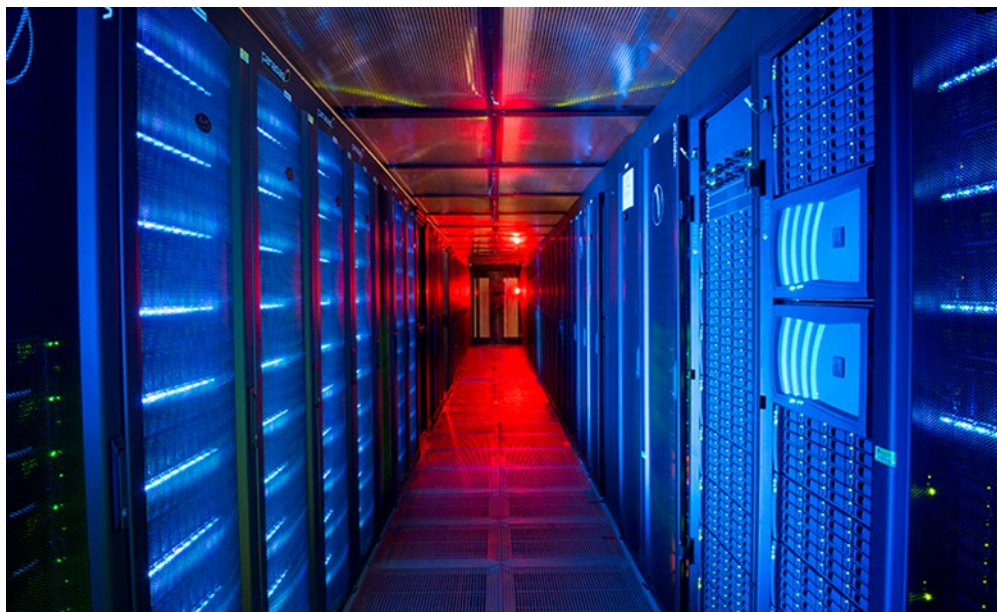# JASMIN-CEMS: Big Data and Compute for Environmental Science

*Victoria Bennett[1,3], Philip Kershaw[1,3], Matt Pritchard[1], Jonathan Churchill[2], Cristina Del Cano Novales[2], Martin Juckes[1,4], Stephen Pascoe[1,4], Sam Pepler[1,4], Ag Stephens[1,4], Bryan Lawrence[1,4,6]*

Centre for Environmental Data Archival, RAL Space, STFC Rutherford Appleton Laboratory, UK; 2. Scientific Computing Department, STFC Rutherford Appleton Laboratory, UK; 3. National Centre for Earth Observation, UK; 4. National Centre for Atmospheric Science, UK; 5. Remote Sensing Group, RAL Space, STFC Rutherford Appleton Laboratory, UK; 6. University of Reading, UK
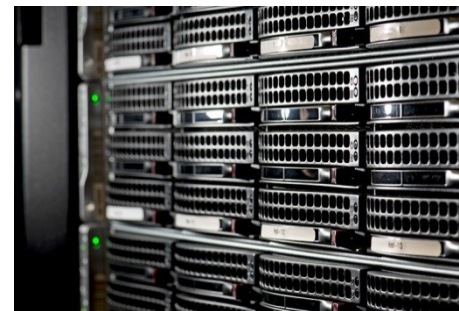
Victoria Bennett
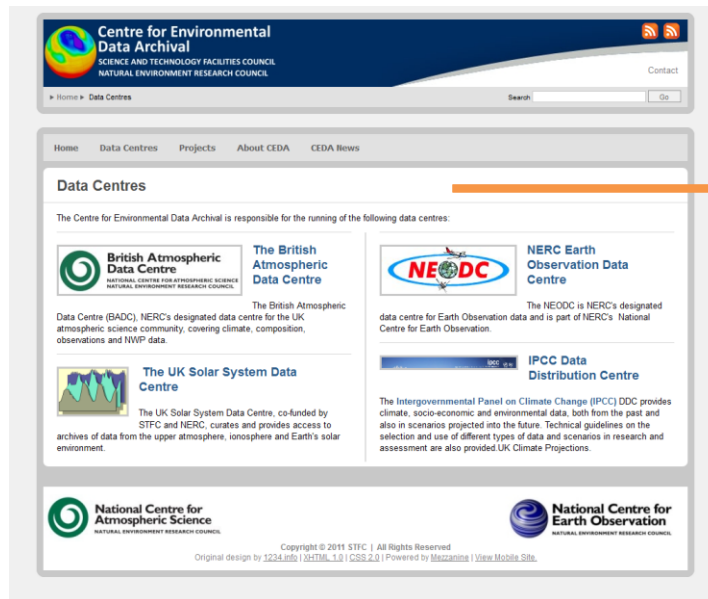
**CEDA, Centre for Environmental Data Archival, STFC**

- Some background
- What is JASMIN, and CEMS
- Facts and figures
- JASMIN operations, and evolution
- Two example science projects



**CEMS**

Climate, Environment &
Monitoring from Space





National Centre for
Atmospheric Science
NATURAL ENVIRONMENT RESEARCH COUNCIL

Centre for Environmental
Data Archival
SCIENCE AND TECHNOLOGY FACILITIES COUNCIL
NATURAL ENVIRONMENT RESEARCH COUNCIL

National Centre for
Earth Observation
NATURAL ENVIRONMENT RESEARCH COUNCIL

# Background: CEDA

## Centre for Environmental Data Archival



"to support environmental science, further environmental data archival practices, and develop and deploy new technologies to enhance access to data"
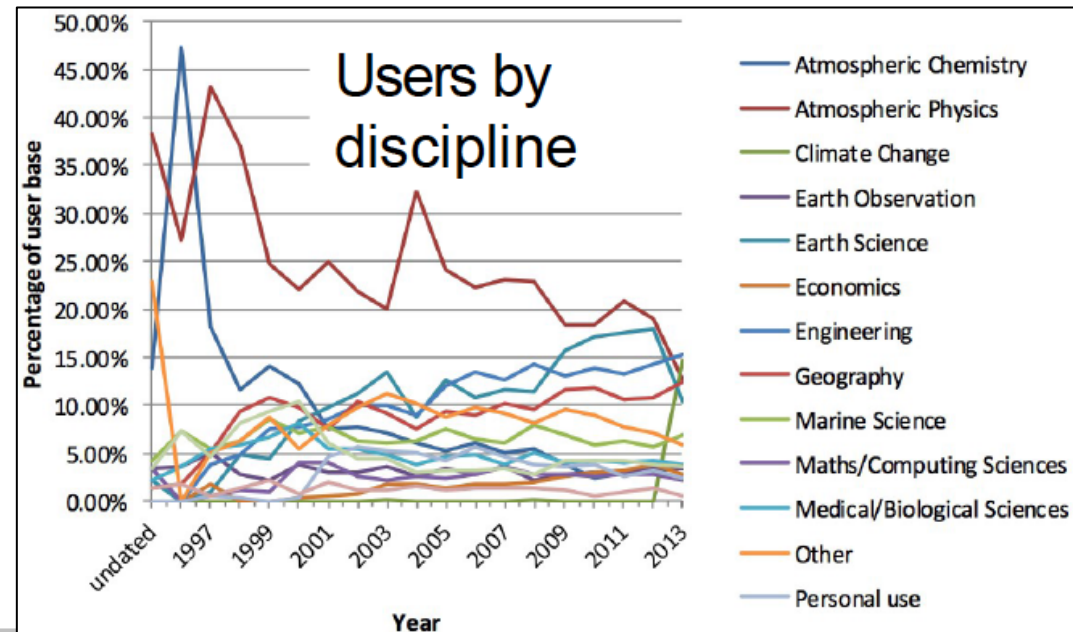
➔ Curation & Facilitation

# Centre for Environmental Data Archival

| Project | Type | Data Volume (Petabytes) |
|---------|------|-------------------------|
| NEODC | Earth Observation | 0.9 |
| BADC | Atmospheric Science | 0.8 |
| CMIP5 | Climate Model | 1.2 |
| | Total | 2.9 |

- > 300 datasets
- 144 million files
- 23,000 registered users



Users by discipline

- Atmospheric Chemistry
- Atmospheric Physics
- Climate Change
- Earth Observation
- Earth Science
- Economics
- Engineering
- Geography
- Marine Science
- Maths/Computing Sciences
- Medical/Biological Sciences
- Other
- Personal use

Growth of Selected Datasets at STFC

(Credit: Folkes, Churchill)

30-85 PB (unique data)
Projection for JASMIN

X?

CMIP6
30-300
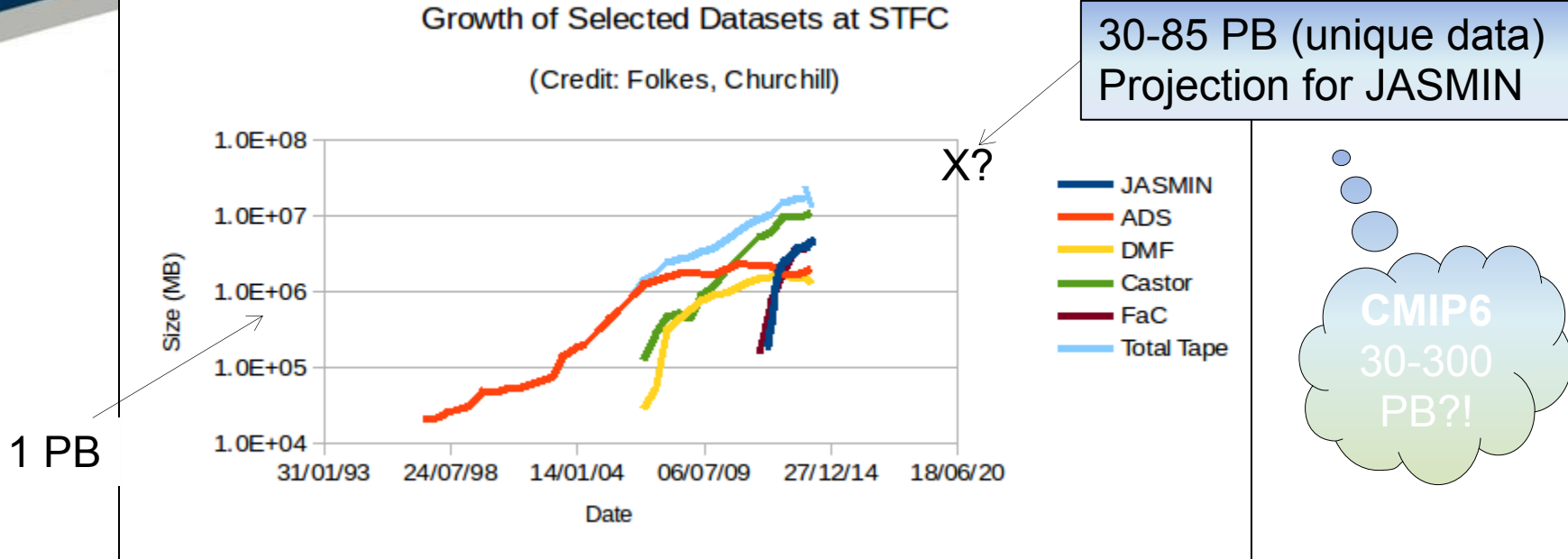PB?!

1 PB
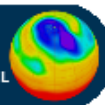
The light blue line is the total of all the data stored on tape in SCD.
The green line is the LHC Tier 1 data on tape.
The dark blue line is the data stored on **disk** in JASMIN.

# JASMIN & CEMS: Big Data Facilities



- JASMIN (super data cluster)
    - Storage and services
    - Scientific computation
    - Access to high volume and complex data
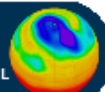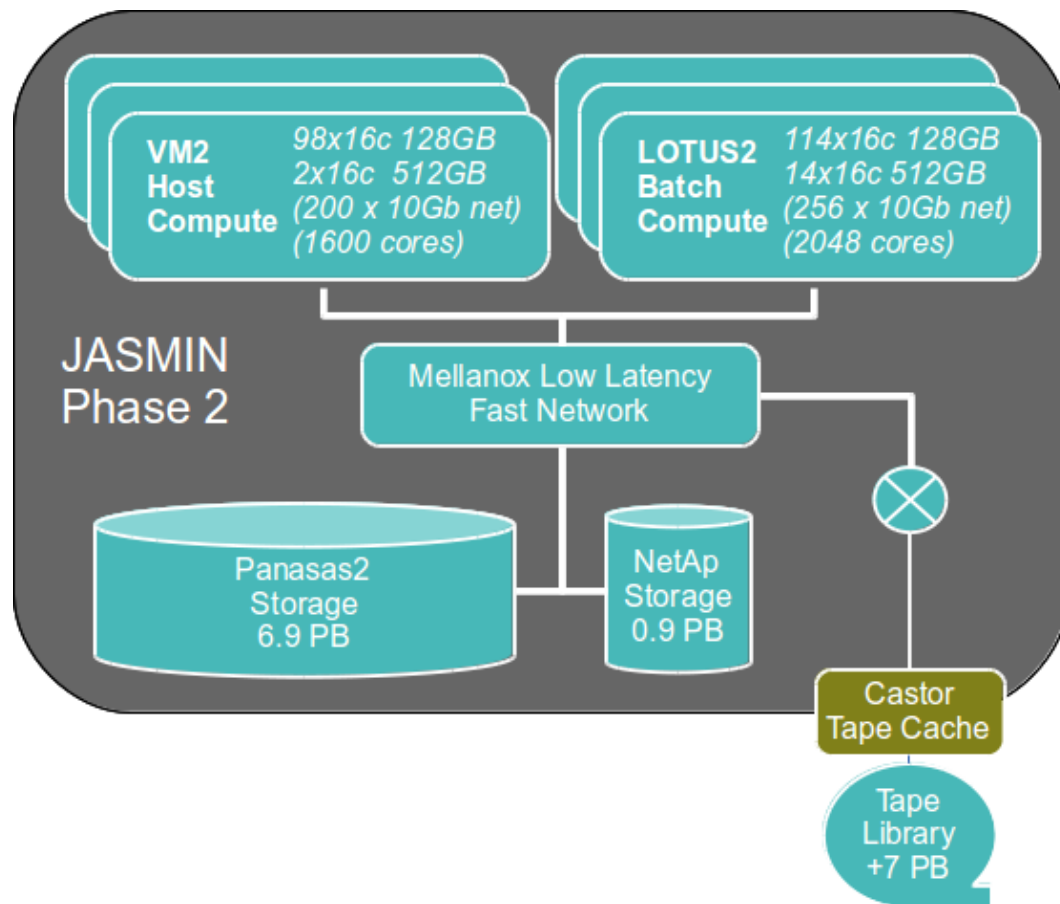
- CEMS: Climate, Environment and Monitoring from Space
    - EO data and services
    - Academic – commercial partnership

# JASMIN

- What have we got:
  - ~16.4 PB fast parallel disk storage & equivalent in near-line tape
  - > 4,000 compute cores

- Four services provided to the community:
  - Storage (disk and tape)
  - Batch computing ("Lotus")
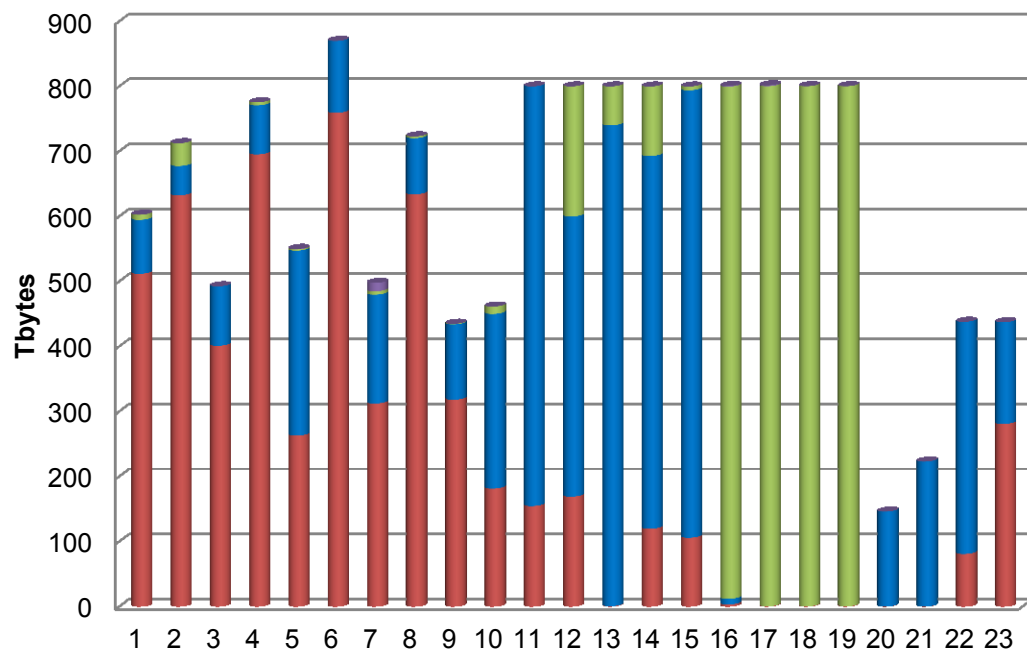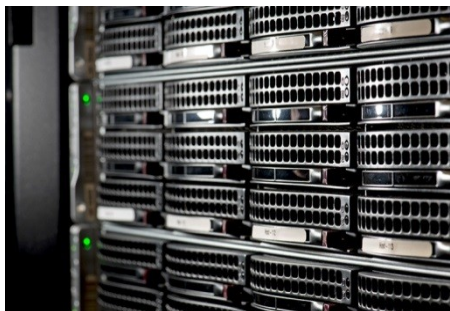  - Hosted computing
  - Cloud computing

- **Data collaboration** – share, process and disseminate;

- **Processing** – analyse own or third-party data;

- **Running models** – port, develop, share and run models;

- **Data/modelling services** – explore, develop and deploy services to provide new interfaces to end-users;

- **Cloud tools** – access to tools that allow creation of virtual servers and allocation of storage resources – for novel research/applications/tools

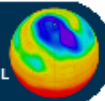**Activities with a research focus**

# JASMIN Operations

- ~600 JASMIN users
- 90 projects
- 5.2 PB allocated as Group Workspace; 3 PB CEDA archives
- Over 2 million processing jobs



JASMIN "bladesets" usage October 2014.
Blue: allocated but not yet used.
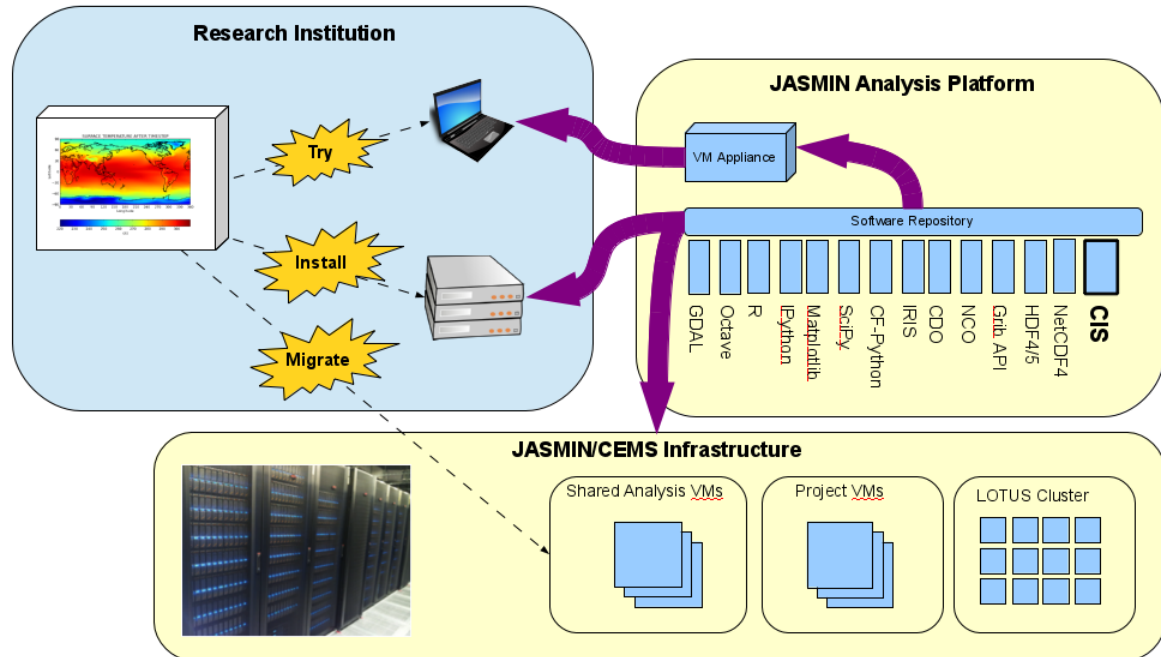Red: used.
Green: as yet unallocated

# JASMIN Analysis Platform (JAP)

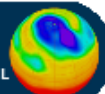Multi-node infrastructure requires a way to install tools quickly and consistently

The community needs a consistent platform where ever they need them.

Users need help migrating analysis to JASMIN.

JAP provides RPMs and pre-built images based on CentOS



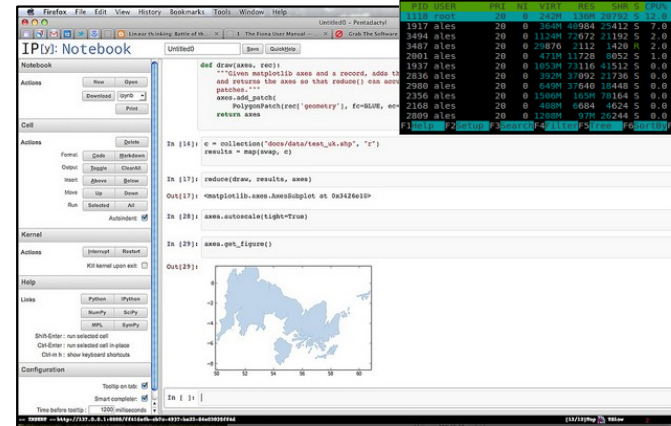http://proj.badc.rl.ac.uk/cedaservices/wiki/JASMIN/AnalysisPlatform

Phase 1:

Storage and batch compute – excellent results for first users..

But "long tail" of user community who are less expert users of e.g. the Linux command line and high performance computing

-> New cloud services to support much wider community
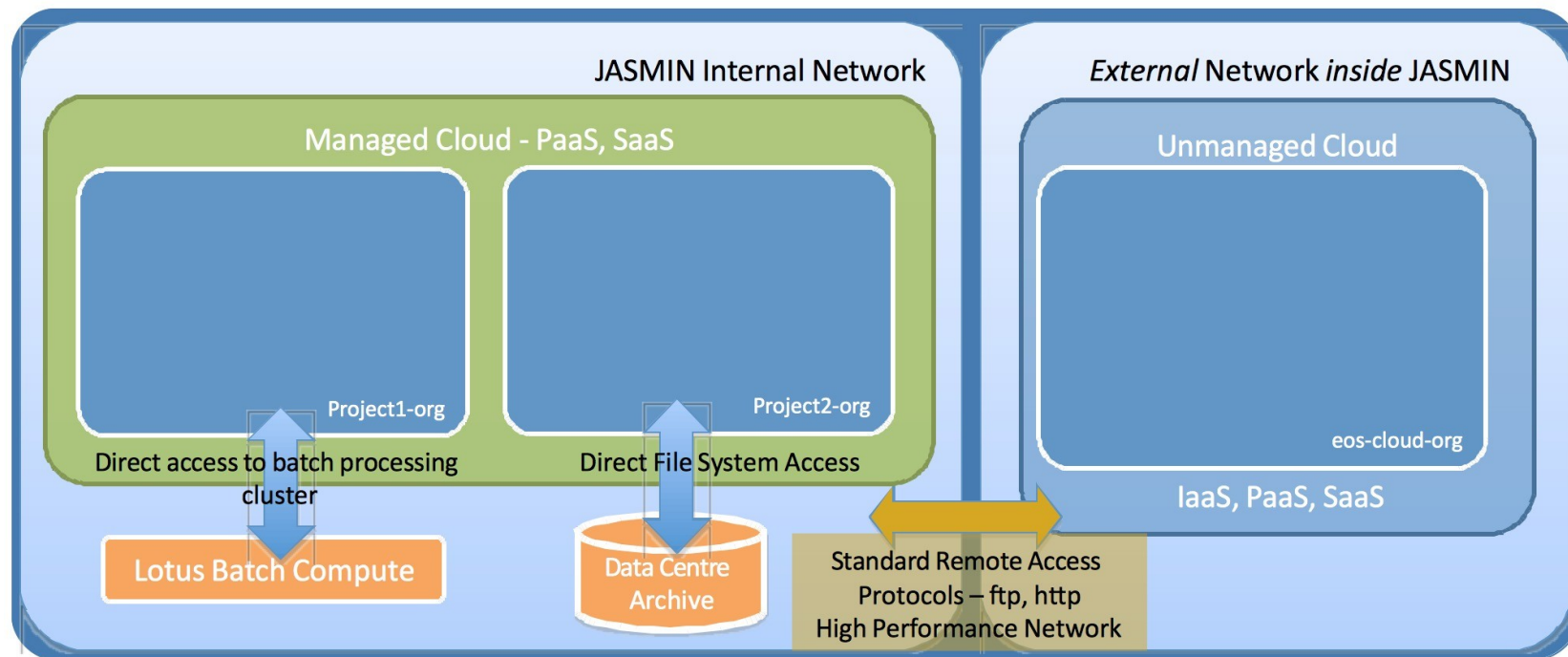


ssh via public IP

IPython Notebook VM could access cluster through Python API

CloudBioLinux Desktop
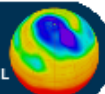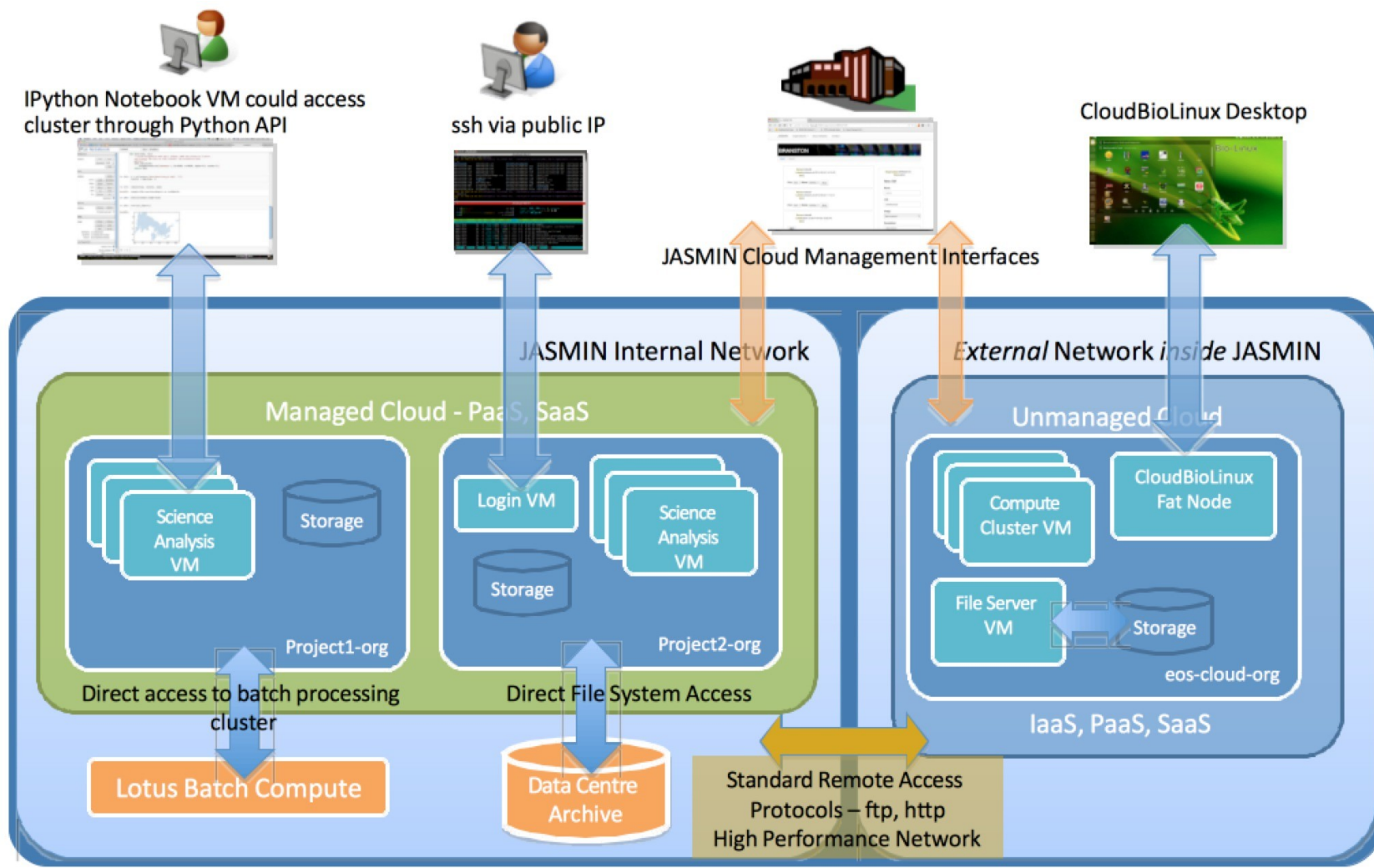
# JASMIN Cloud Architecture



Group Work Spaces and hosted processing in the Managed Cloud:
direct access to archive filesystem and Lotus batch processing
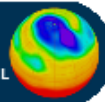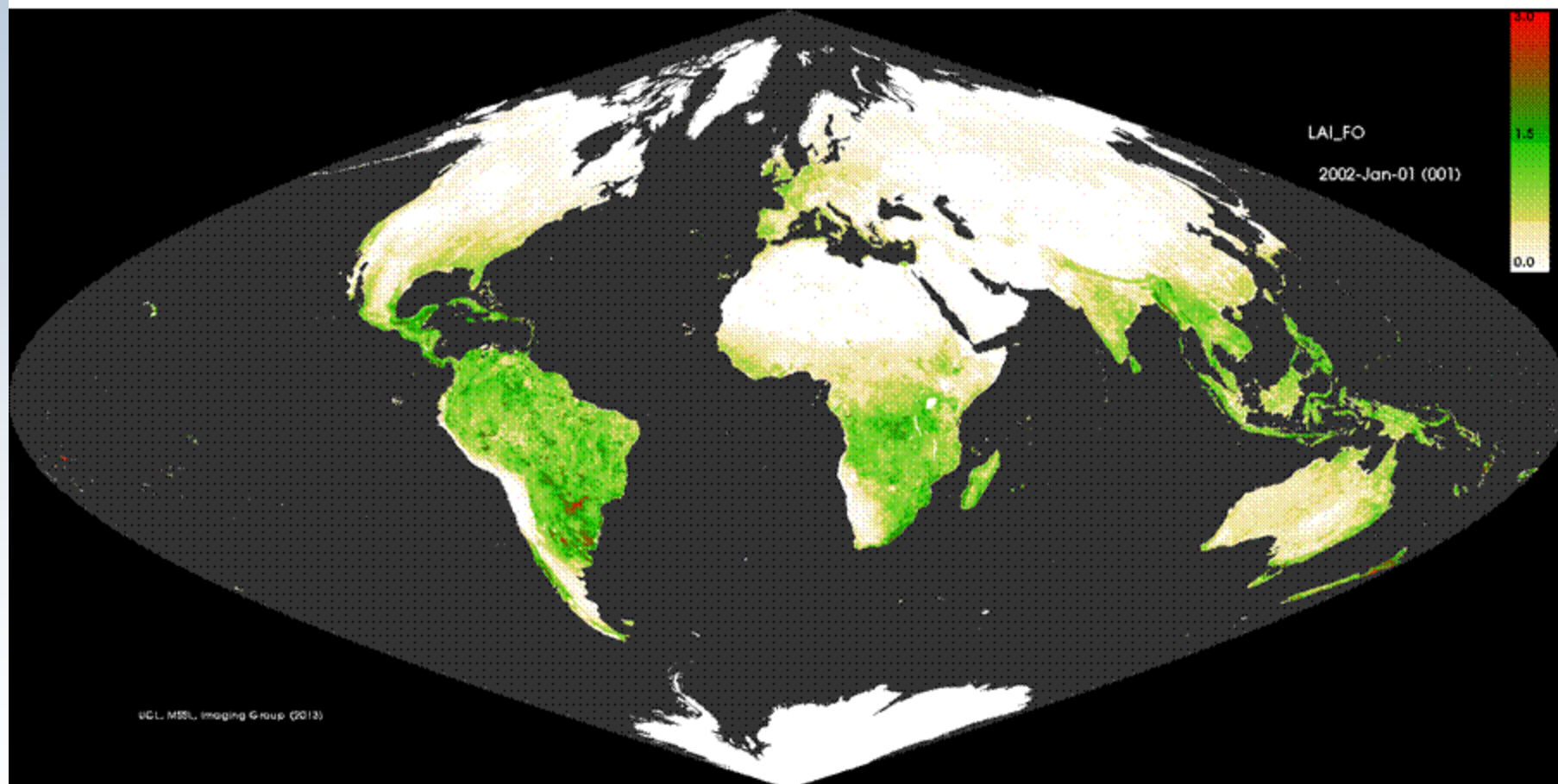
First projects underway in the Unmanaged Cloud

# Example EO Science projects using JASMIN-CEMS

# JASMIN-CEMS for global land surface products

- **Objective 1:** Re-project BRDF files from SIN-coordinates to lat/lon
  - **Challenge**: huge number of polygons to be spatiality indexed and processed. **This process requires massive RAM and usually takes a very long time!**

- **Objective 2:** Create specific albedo products for computation of 8-daily LAI/fAPAR between 2002 and 2011 at 3 different resolutions: 1km, 5km and 25km
  - **Challenge:** Upscale big data BRDF (50TB) from 1km to 5km and 25km using energy conservation method: **This process is extremely time consuming!**

- **Solution**: Cloud-computing system in JASMIN-CEMS (~100 times faster  than 224-core in house linux cluster)

- Also use Science DMZ for data transfers from NASA: achieved rates up to 28 TB/day
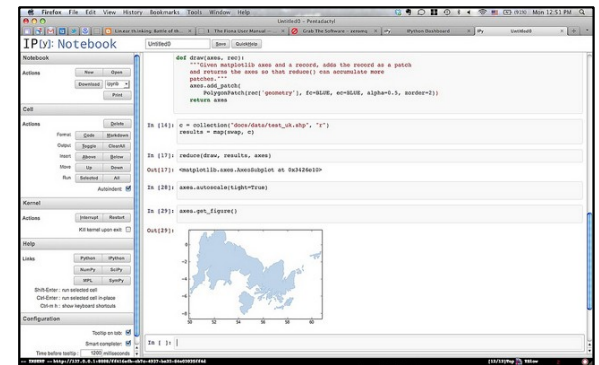
# Example of 8-daily Global LAI derived from GlobAlbedo for 10 years using combined processing on CEMS and FastOpt Hamburg

# ESA OPTIRAD Project

- Developing a Collaborative Research Environment for land data assimilation
  - a dedicated **software environment for the scientific community to generate products** from raw EO data
  - compute intensive assimilation algorithms with high memory demands
- Using iPython Notebook  on the CEMS Unmanaged cloud

# Further info

JASMIN

http://www.jasmin.ac.uk

Centre for Environmental Data Archival

http://www.ceda.ac.uk

JASMIN paper

Lawrence, B.N. , V.L. Bennett, J. Churchill, M. Juckes, P. Kershaw, S. Pascoe, S. Pepler, M. Pritchard, and A. Stephens. **Storing and manipulating environmental big data with JASMIN.** *Proceedings of IEEE Big Data 2013, p68-75,* doi:10.1109/BigData.2013.6691556

victoria.bennett@stfc.ac.uk